

Decision Trees

Nipun Batra and teaching staff

IIT Gandhinagar

August 1, 2025

Table of Contents

1. Introduction and Motivation
2. Discrete Input, Real Output
3. Real Input Discrete Output
4. Real Input Real Output
5. Pruning and Overfitting
6. Summary and Key Takeaways
7. Weighted Entropy

The need for interpretability

How to maintain trust in AI

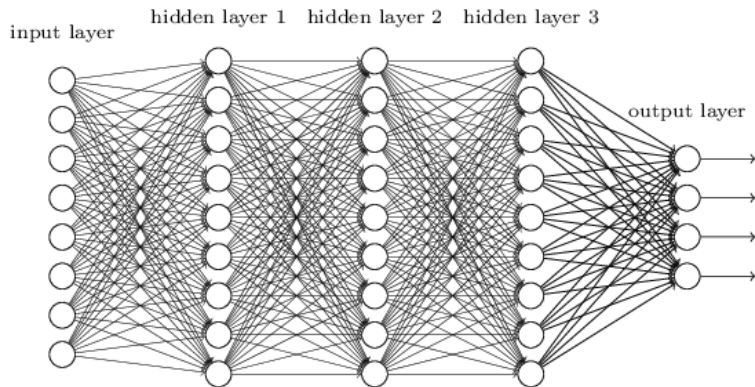
Beyond developing initial trust, however, creators of AI also must work to maintain that trust. Siau and Wang suggest seven ways of "developing continuous trust" beyond the initial phases of product development:

- Usability and reliability. AI "should be designed to operate easily and intuitively," Siau and Wang write. "There should be no unexpected downtime or crashes."
- Collaboration and communication. AI developers want to create systems that perform autonomously, without human involvement. Developers must focus on creating AI applications that smoothly and easily collaborate and communicate with humans.
- Sociability and bonding. Building social activities into AI applications is one way to strengthen trust. A robotic dog that can recognize its owner and show affection is one example, Siau and Wang write.
- Security and privacy protection. AI applications rely on large data sets, so ensuring privacy and security will be crucial to establishing trust in the applications.
- Interpretability. Just as transparency is instrumental in building initial trust, interpretability – or the ability for a machine to explain its conclusions or actions – will help sustain trust.

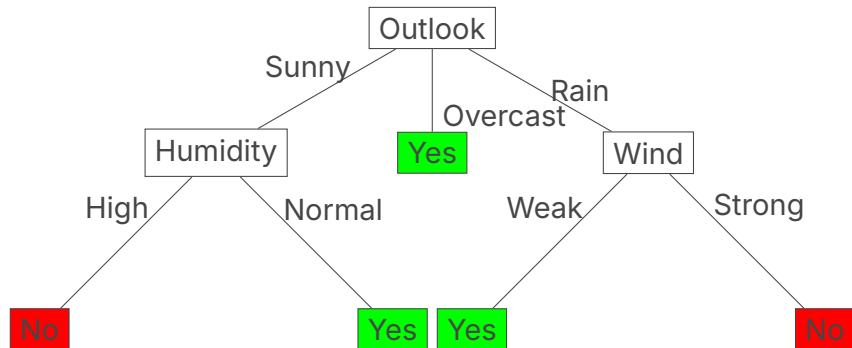
Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Learning a Complicated Neural Network



Learnt Decision Tree



Medical Diagnosis using Decision Trees



Source: Improving medical decision trees by combining relevant health-care criteria

Leo Breiman



Leo Breiman 1928-2005

Professor of Statistics, UC Berkeley
Verified email at stat.berkeley.edu - [Homepage](#)
[Data Analysis](#) [Statistics](#) [Machine Learning](#)

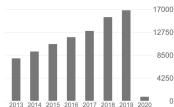
[FOLLOW](#)

TITLE	CITED BY	YEAR
Random forests L. Breiman Machine learning 45 (1), 5-32	53816	2001
Classification and Regression Trees L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone CRC Press, New York	43992 *	1999
Classification and regression trees L. Breiman Chapman & Hall/CRC	43992 *	1984
Bagging predictors L. Breiman Machine learning 24 (2), 123-140	22742	1996
Statistical Modeling: The Two Cultures L. Breiman	2788 *	2003
Statistical modeling: The two cultures (with comments and a rejoinder by the author) L. Breiman Statistical Science 16 (3), 199-231	2772	2001
Estimating optimal transformations for multiple regression and correlation	2096	1985

Cited by

[VIEW ALL](#)

	All	Since 2015
Citations	142857	68736
h-index	51	33
i10-index	80	46



Optimal Decision Tree

Volume 5, number 1

INFORMATION PROCESSING LETTERS

May 1976

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL

IRIA – Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST

Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Received 7 November 1975, revised version received 26 January 1976

Binary decision trees, computational complexity, NP-complete

Pop Quiz #1

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features

Pop Quiz #2

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node

Pop Quiz #3

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data

Pop Quiz #4

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Pop Quiz #5

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Pop Quiz #6

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Answer: D) All of the above - The search space is exponentially large, making brute force optimization computationally intractable.

Greedy Algorithm

Core idea: At each level, choose an attribute that gives **biggest estimated** performance gain!

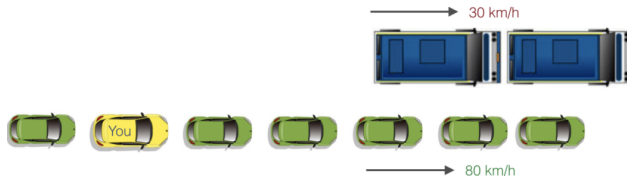


Image source: analyticsvidhya

Greedy \neq Optimal

Towards biggest estimated performance gain

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is “easier” when there is less disagreement

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is “easier” when there is less disagreement

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is "easier" when there is less disagreement
- Need some statistical measure of "disagreement"

Entropy

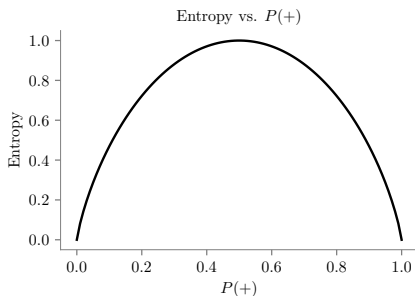
Statistical measure to characterize the (im)purity of examples

Entropy

Statistical measure to characterize the (im)purity of examples

$$H(X) = - \sum_{i=1}^k p(x_i) \log_2 p(x_i)$$

Notebook: entropy.html



Towards biggest estimated performance gain

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?
- When Outlook is overcast, we always Play and thus no "disagreement"

Information Gain

Reduction in entropy by partitioning examples (S) on attribute A

$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Pop Quiz #7

Quick Question!

What does entropy measure in the context of decision trees?

A) The depth of the tree

Pop Quiz #8

Quick Question!

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples

Pop Quiz #9

Quick Question!

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset

Pop Quiz #10

Quick Question!

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

Pop Quiz #11

Quick Question!

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

Pop Quiz #12

Quick Question!

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

Answer: B) The impurity or "disagreement" in a set of examples - Higher entropy means more mixed classes, lower entropy means more pure subsets.

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - Examples_v : subset of examples that $A = v$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - Examples_v : subset of examples that $A = v$
 - If Examples_v is empty: add leaf with label = most common value of Target Attribute

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $Root \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - $Examples_v$: subset of examples that $A = v$
 - If $Examples_v$ is empty: add leaf with label = most common value of Target Attribute
 - Else: ID3 ($Examples_v$, Target attribute, Attributes - A)

Learnt Decision Tree

Root Node (empty)

Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Entropy calculated

We have 14 examples in S : 5 No, 9 Yes

$$\begin{aligned}\text{Entropy}(S) &= -p_{\text{No}} \log_2 p_{\text{No}} - p_{\text{Yes}} \log_2 p_{\text{Yes}} \\ &= -\frac{5}{14} \log_2 \left(\frac{5}{14} \right) - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) = 0.940\end{aligned}$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Overcast	Yes
Rain	Yes
Rain	Yes
Rain	No
Overcast	Yes
Sunny	No
Sunny	Yes
Rain	Yes
Sunny	Yes
Overcast	Yes
Overcast	Yes
Rain	No

Information Gain for Outlook

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3

$$\begin{aligned} \text{No Entropy} = & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3

$$\begin{aligned} \text{No Entropy} = & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3

$$\begin{aligned} \text{No Entropy} = & \\ & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

Outlook	Play
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

We have 4 Yes, 0

No Entropy = 0
(pure subset)

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3

$$\begin{aligned} \text{No Entropy} = & \\ & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

Outlook	Play
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

We have 4 Yes, 0

No Entropy = 0
(pure subset)

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3

$$\begin{aligned} \text{No Entropy} = & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

Outlook	Play
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

We have 4 Yes, 0

$$\begin{aligned} \text{No Entropy} = & 0 \\ & (\text{pure subset}) \end{aligned}$$

Outlook	Play
Rain	Yes
Rain	Yes
Rain	No
Rain	Yes
Rain	No

We have 3 Yes, 2

$$\begin{aligned} \text{No Entropy} = & -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \\ & \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \end{aligned}$$

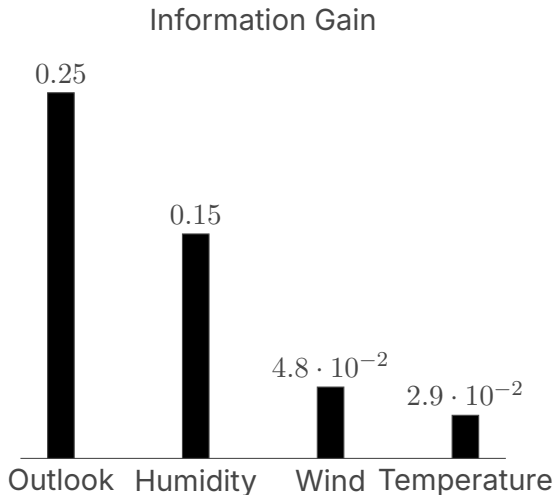
Information Gain

$$\text{Gain}(\mathcal{S}, \text{Outlook}) = \text{Entropy}(\mathcal{S}) - \sum_{v \in \{\text{Rain, Sunny, Overcast}\}} \frac{|\mathcal{S}_v|}{|\mathcal{S}|} \text{Entropy}(\mathcal{S}_v)$$

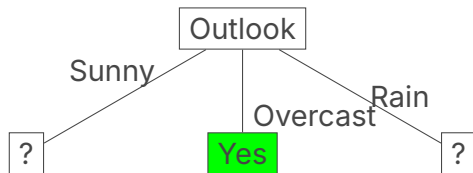
$$\text{Gain}(\mathcal{S}, \text{Outlook}) = \text{Entropy}(\mathcal{S}) - \frac{5}{14} \text{Entropy}(\mathcal{S}_{\text{Sunny}}) - \frac{4}{14} \text{Entropy}(\mathcal{S}_{\text{Overcast}})$$

$$= 0.940 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.971 = 0.940 - 0.347 - 0 - 0.347 = 0.246$$

Information Gain



Learnt Decision Tree



Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain($S_{\text{Outlook=Sunny}}$, Temp) = Entropy(2 Yes, 3 No) -
(2/5)*Entropy(0 Yes, 2 No) - (2/5)*Entropy(1 Yes, 1 No)
- (1/5)*Entropy(1 Yes, 0 No)

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain($S_{\text{Outlook=Sunny}}$, Temp) = Entropy(2 Yes, 3 No) -
(2/5)*Entropy(0 Yes, 2 No) - (2/5)*Entropy(1 Yes, 1 No)
- (1/5)*Entropy(1 Yes, 0 No)

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Temp}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(0 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(1 \text{ Yes}, 1 \text{ No}) - (1/5) * \text{Entropy}(1 \text{ Yes}, 0 \text{ No})$
- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Humidity}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ Yes}, 0 \text{ No}) - (3/5) * \text{Entropy}(0 \text{ Yes}, 3 \text{ No}) \Rightarrow$ **maximum possible for the set**

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

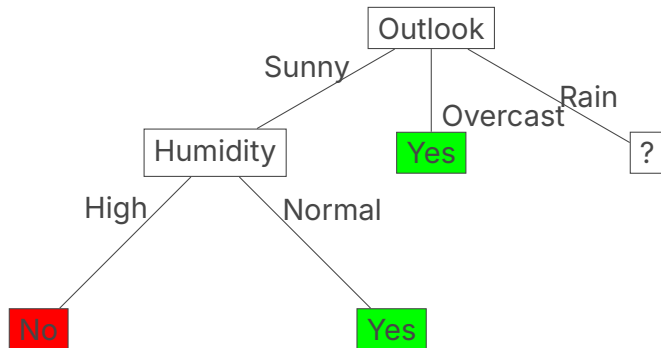
- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Temp}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(0 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(1 \text{ Yes}, 1 \text{ No}) - (1/5) * \text{Entropy}(1 \text{ Yes}, 0 \text{ No})$
- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Humidity}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ Yes}, 0 \text{ No}) - (3/5) * \text{Entropy}(0 \text{ Yes}, 3 \text{ No}) \Rightarrow$ **maximum possible for the set**

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Temp}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(0 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(1 \text{ Yes}, 1 \text{ No}) - (1/5) * \text{Entropy}(1 \text{ Yes}, 0 \text{ No})$
- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Humidity}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ Yes}, 0 \text{ No}) - (3/5) * \text{Entropy}(0 \text{ Yes}, 3 \text{ No}) \implies \textbf{maximum possible for the set}$
- $\text{Gain}(S_{\text{Outlook=Sunny}}, \text{Windy}) = \text{Entropy}(2 \text{ Yes}, 3 \text{ No}) - (3/5) * \text{Entropy}(1 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(1 \text{ Yes}, 1 \text{ No})$

Learnt Decision Tree

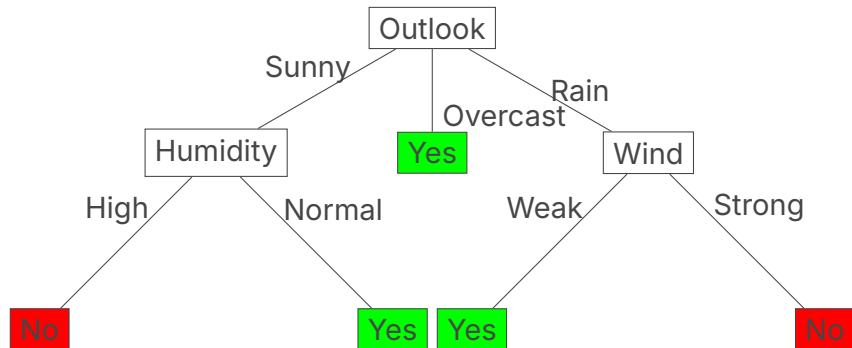


Calling ID3 on (Outlook=Rain)

Day	Temp	Humidity	Windy	Play
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

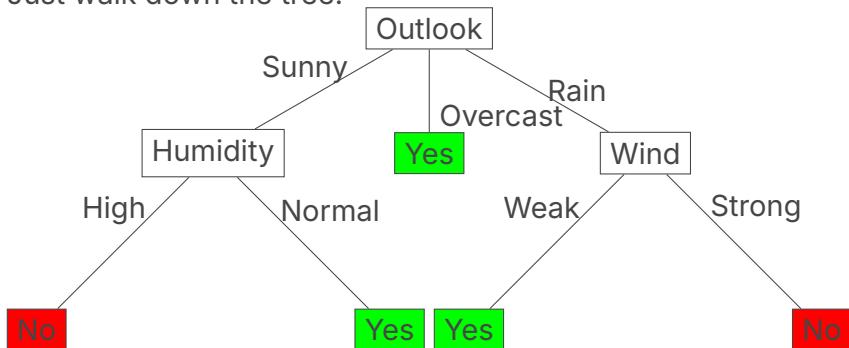
- The attribute Windy gives the highest information gain

Learnt Decision Tree



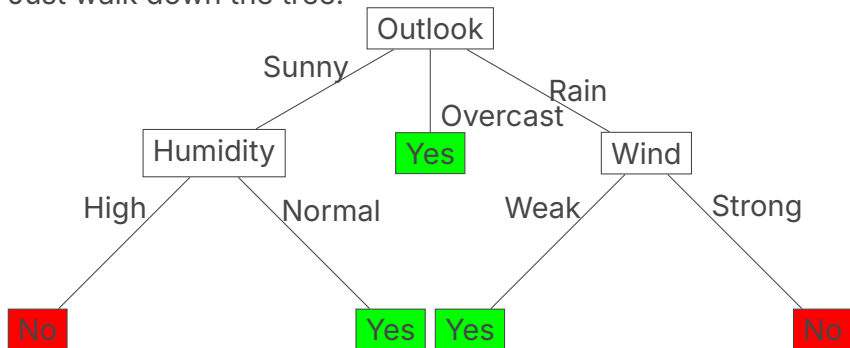
Prediction for Decision Tree

Just walk down the tree!



Prediction for Decision Tree

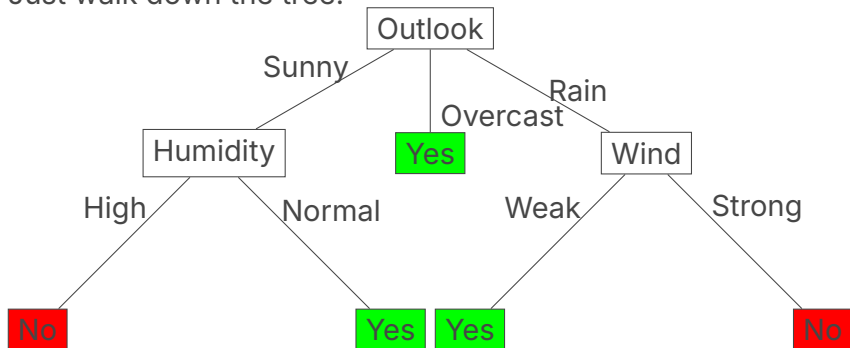
Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ?

Prediction for Decision Tree

Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ?
No

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

What is depth-1 tree (no decision) for the examples?

Limiting Depth of Tree

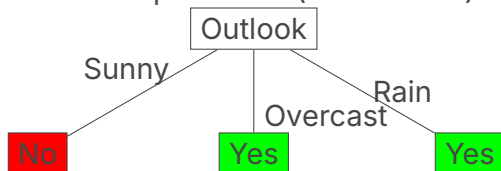
Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

What is depth-1 tree (no decision) for the examples?



Pop Quiz #13

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

A) It was the first feature in the dataset

Pop Quiz #14

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)

Pop Quiz #15

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values

Pop Quiz #16

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Pop Quiz #17

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Pop Quiz #18

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Answer: B) When Outlook=Overcast, all examples have Play=Yes - This creates a pure subset with entropy=0, maximizing information gain.

Modified Dataset

Day	Outlook	Temp	Humidity	Wind	Minutes Played
D1	Sunny	Hot	High	Weak	20
D2	Sunny	Hot	High	Strong	24
D3	Overcast	Hot	High	Weak	40
D4	Rain	Mild	High	Weak	50
D5	Rain	Cool	Normal	Weak	60
D6	Rain	Cool	Normal	Strong	10
D7	Overcast	Cool	Normal	Strong	4
D8	Sunny	Mild	High	Weak	10
D9	Sunny	Cool	Normal	Weak	60
D10	Rain	Mild	Normal	Weak	40
D11	Sunny	Mild	High	Strong	45
D12	Overcast	Mild	High	Strong	40
D13	Overcast	Hot	Normal	Weak	35
D14	Rain	Mild	High	Strong	20

Measure of Impurity for Regression?

Measure of Impurity for Regression?

Measure of Impurity for Regression?

- Any guesses?

Measure of Impurity for Regression?

- Any guesses?

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$
- What about splitting criterion for regression?

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$
- What about splitting criterion for regression?

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$
- What about splitting criterion for regression?
- **MSE Reduction** (not Information Gain!)

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$
- What about splitting criterion for regression?
- **MSE Reduction** (not Information Gain!)

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$
- What about splitting criterion for regression?
- **MSE Reduction** (not Information Gain!)
- $\text{MSE Reduction} = \text{MSE}(S) - \sum_v \frac{|S_v|}{|S|} \text{MSE}(S_v)$

Gain by splitting on Wind

Gain by splitting on Wind

Wind	Minutes Played
Weak	20
Strong	24
Weak	40
Weak	50
Weak	60
Strong	10
Strong	4
Weak	10
Weak	60
Weak	40
Strong	45
Strong	40
Weak	35
Strong	20

$$\text{MSE}(S)=311.34$$

Gain by splitting on Wind

Wind	Minutes Played
Weak	20
Strong	24
Weak	40
Weak	50
Weak	60
Strong	10
Strong	4
Weak	10
Weak	60
Weak	40
Strong	45
Strong	40
Weak	35
Strong	20

$$\text{MSE}(S) = 311.34$$

Wind	Minutes Played
Weak	20
Weak	40
Weak	50
Weak	60
Weak	10
Weak	60
Weak	40
Weak	35

$$\text{MSE}(S_{\text{Wind=Weak}}) = 277, \text{ Weight} = \frac{8}{14}$$

Wind	Minutes Played
Strong	24
Strong	10
Strong	4
Strong	45
Strong	40
Strong	20

$$\text{MSE}(S_{\text{Wind=Strong}}) = 218, \text{ Weight} = \frac{6}{14}$$

MSE Reduction Calculation

Correct calculation for Wind split:

MSE Reduction = $\text{MSE}(S)$ – Weighted Average MSE

$$= 311.34 - \left[\frac{8}{14} \times 277 + \frac{6}{14} \times 218 \right] = 311.34 - [158.857 + 93.429] = 311.34 - 252.286 = 59.054$$

Key insight: MSE Reduction > 0 means the split improves our model!

For regression: Use MSE Reduction, NOT Information Gain!

Pop Quiz #19

Quick Question!

For regression trees, what criterion do we use instead of Information Gain?

A) Information Gain

Pop Quiz #20

Quick Question!

For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity

Pop Quiz #21

Quick Question!

For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction

Pop Quiz #22

Quick Question!

For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

Pop Quiz #23

Quick Question!

For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

Pop Quiz #24

Quick Question!

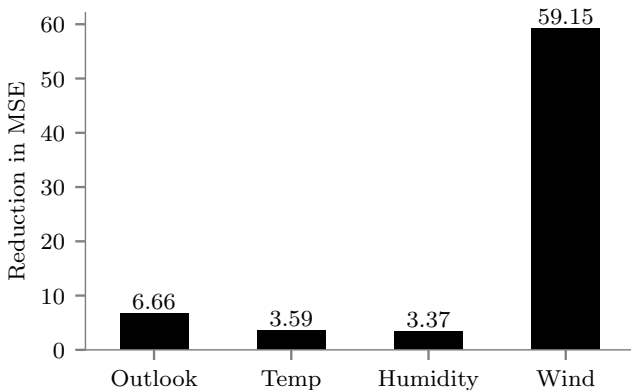
For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

Answer: C) Mean Squared Error (MSE) Reduction
- For regression, we minimize MSE instead of maximizing information gain.

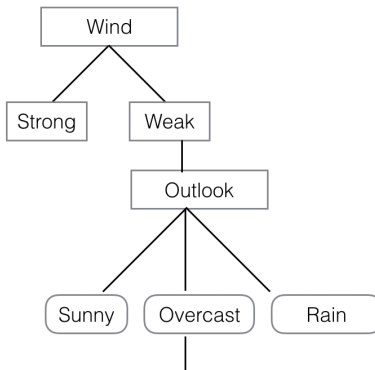
MSE Reduction for Regression Trees

Notebook: decision-tree-real-output.html



Learnt Tree

Assume a tree like this is learnt ...



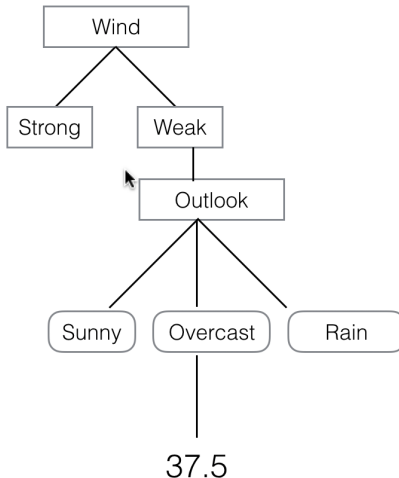
	Day	Outlook	Temp	Humidity	Wind	Minutes Played
2	D3	Overcast	Hot	High	Weak	40
12	D13	Overcast	Hot	Normal	Weak	35

Learnt Tree

Method 1

Mins

Played=(40+35)
/2



Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split
- Choose the split with the lowest impurity

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No
- Entropy for LHS = 0, Entropy for RHS = 0.971

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No
- Entropy for LHS = 0, Entropy for RHS = 0.971
- Weighted Entropy = $0.971 \times 5/6 = 0.808$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No
- Entropy for LHS = 0, Entropy for RHS = 0.811

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No
- Entropy for LHS = 0, Entropy for RHS = 0.811
- Weighted Entropy = $0.811 \times 4/6 = 0.541$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No
- Entropy for LHS = 0.918, Entropy for RHS = 0.918

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No
- Entropy for LHS = 0.918, Entropy for RHS = 0.918
- Weighted Entropy = $0.918 \times 3/6 + 0.918 \times 3/6 = 0.918$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No
- Entropy for LHS = 1, Entropy for RHS = 1

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

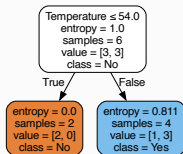
- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No
- Entropy for LHS = 1, Entropy for RHS = 1
- Weighted Entropy = $1 \cdot 4/6 + 1 \cdot 2/6 = 1$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook:
output.html

decision-tree-real-input-discrete-

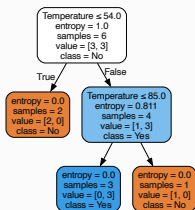


Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook:
output.html

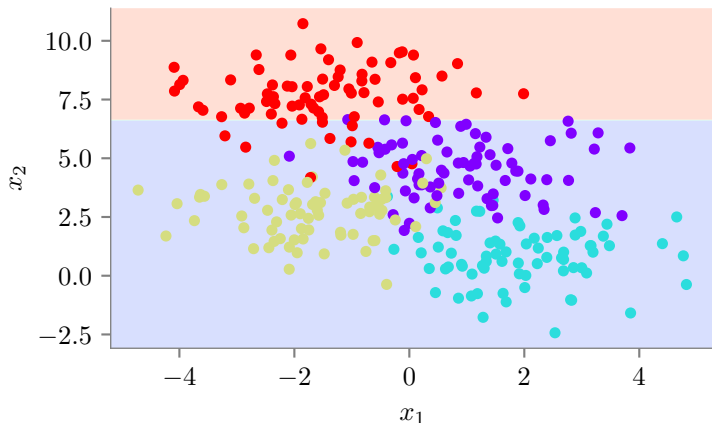
decision-tree-real-input-discrete-



Example (DT of depth 1)

Notebook:
output.html

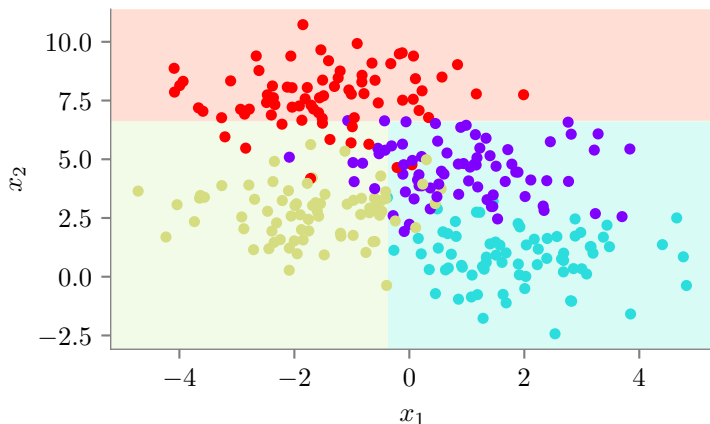
decision-tree-real-input-discrete-



Example (DT of depth 2)

Notebook:
output.html

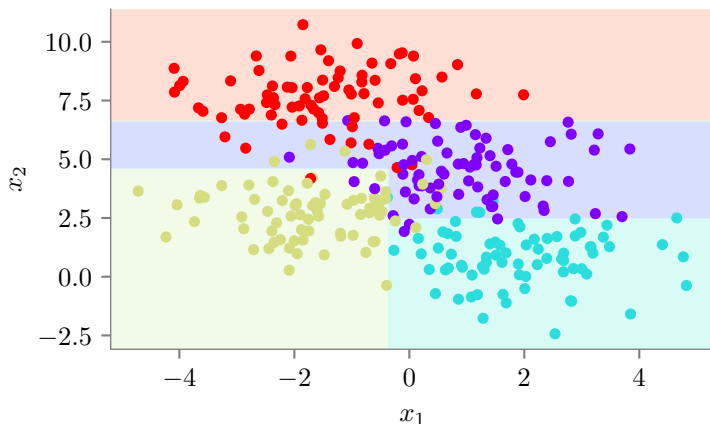
decision-tree-real-input-discrete-



Example (DT of depth 3)

Notebook:
output.html

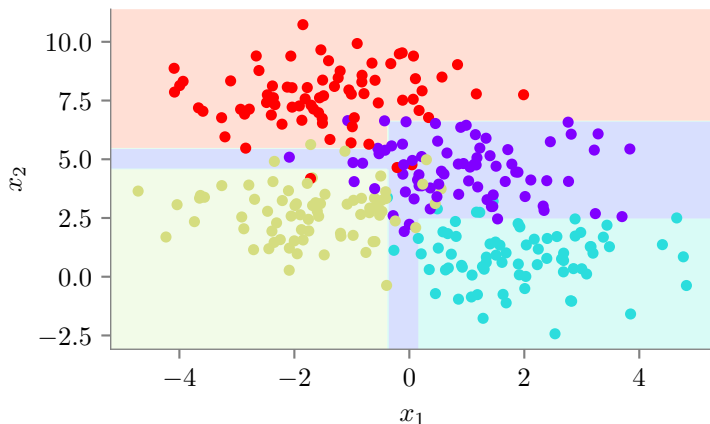
decision-tree-real-input-discrete-



Example (DT of depth 4)

Notebook:
output.html

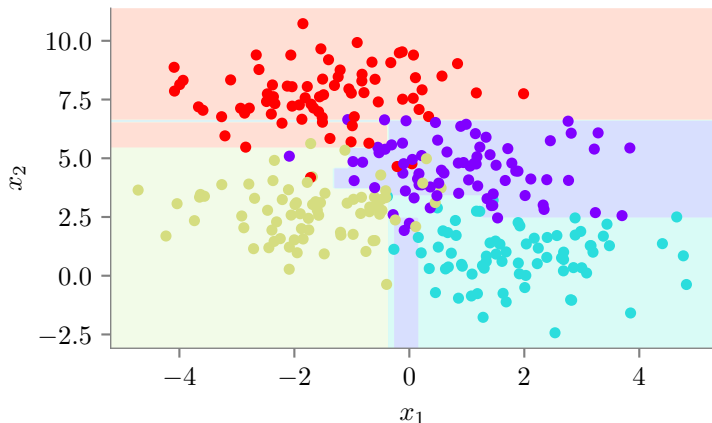
decision-tree-real-input-discrete-



Example (DT of depth 5)

Notebook:
output.html

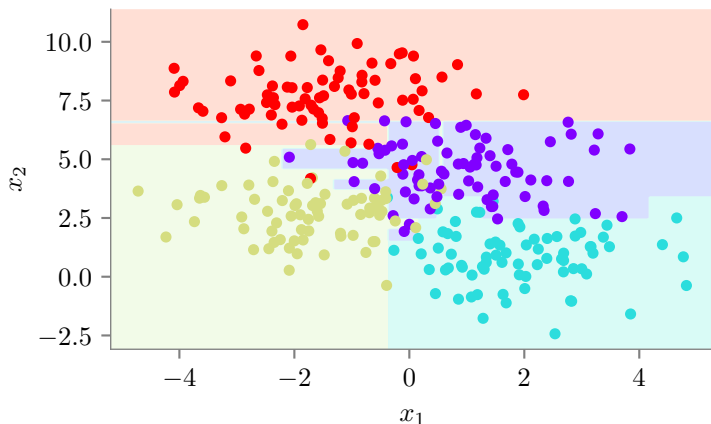
decision-tree-real-input-discrete-



Example (DT of depth 6)

Notebook:
output.html

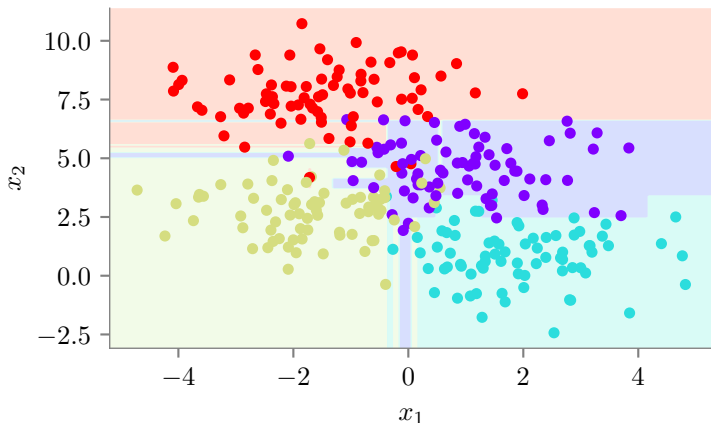
decision-tree-real-input-discrete-



Example (DT of depth 7)

Notebook:
output.html

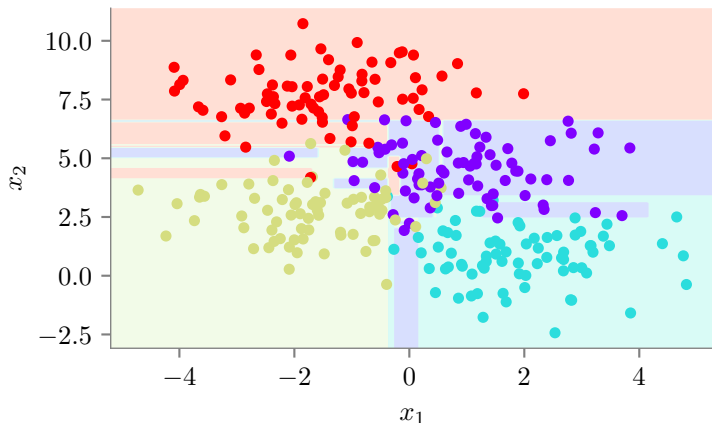
decision-tree-real-input-discrete-



Example (DT of depth 8)

Notebook:
output.html

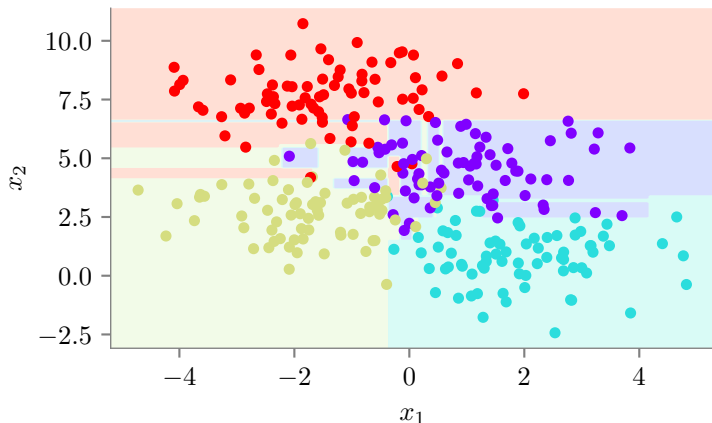
decision-tree-real-input-discrete-



Example (DT of depth 9)

Notebook:
output.html

decision-tree-real-input-discrete-



Pop Quiz #25

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

A) Use all feature values as split points

Pop Quiz #26

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values

Pop Quiz #27

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range

Pop Quiz #28

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Pop Quiz #29

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Pop Quiz #30

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

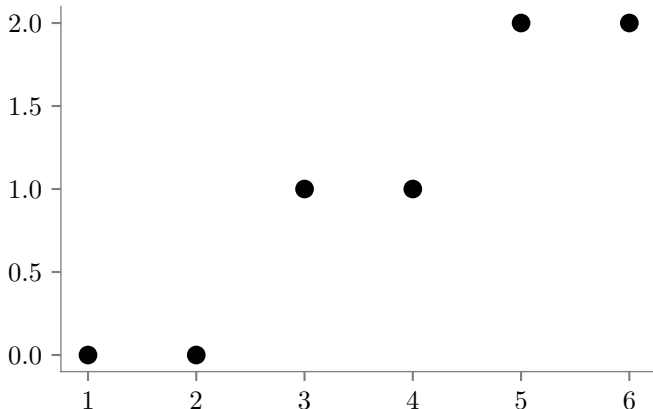
- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Answer: B) Use midpoints between consecutive sorted feature values - This ensures we test all meaningful boundaries between different class regions.

Example 1

Let us consider the dataset given below

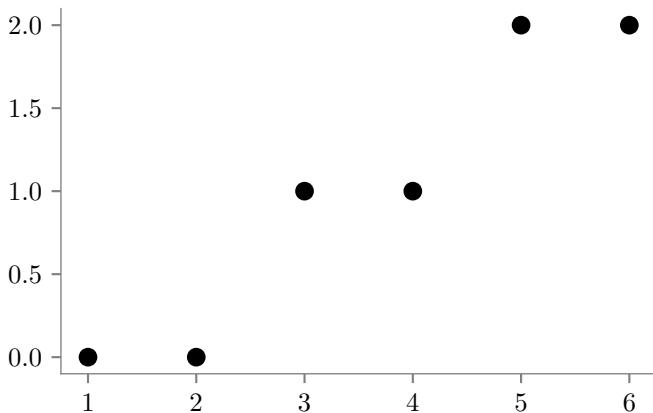
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

What would be the prediction for decision tree with depth 0?

Notebook: [decision-tree-real-input-real-output.html](#)

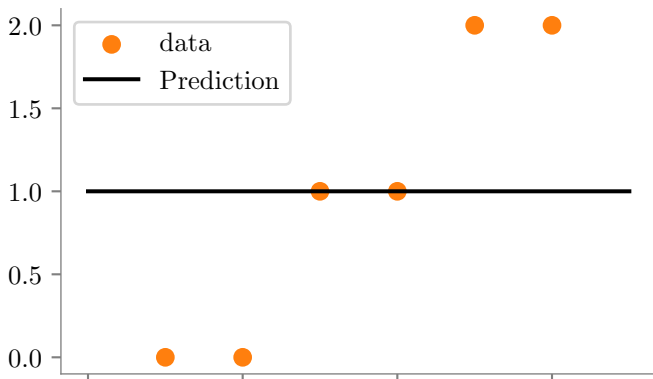


Example 1

Prediction for decision tree with depth 0.

Horizontal dashed line shows the predicted Y value. It is the average of Y values of all datapoints.

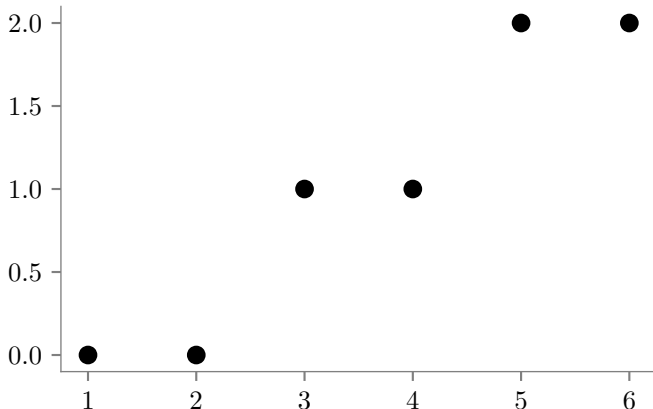
Notebook: decision-tree-real-input-real-output.html



Example 1

What would be the decision tree with depth 1?

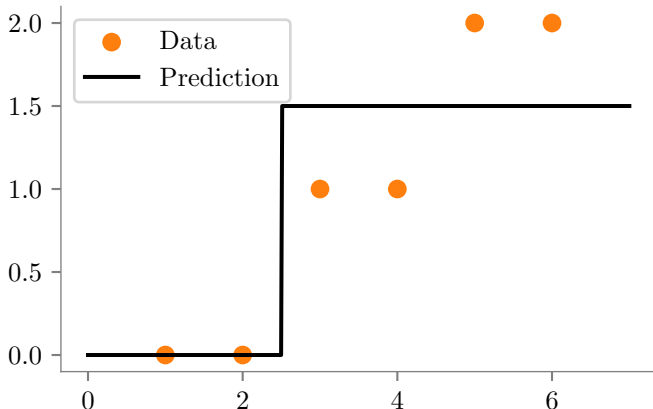
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

Decision tree with depth 1

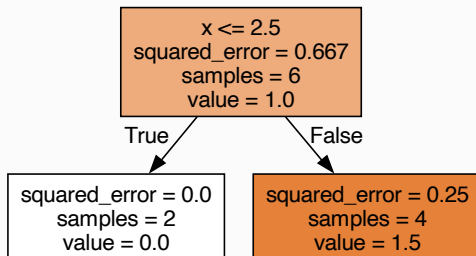
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

The Decision Boundary

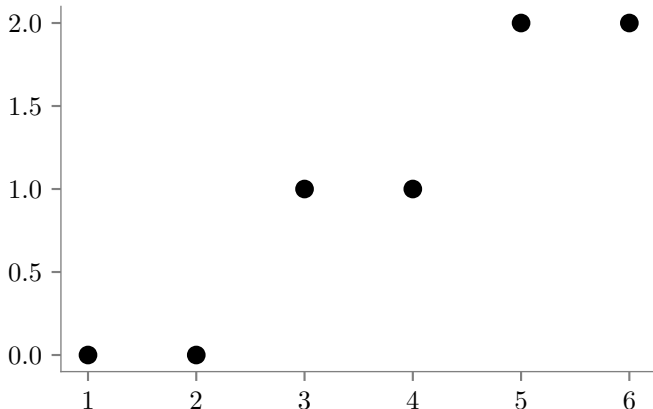
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

What would be the decision tree with depth 2 ?

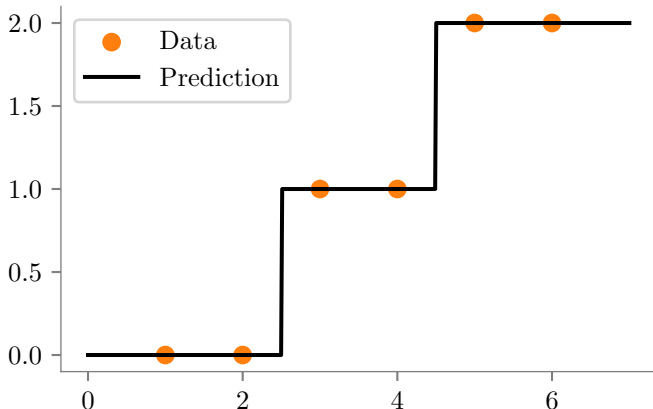
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

Decision tree with depth 2

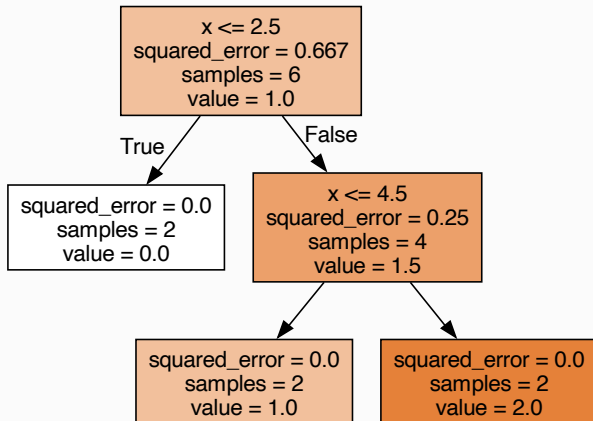
Notebook: [decision-tree-real-input-real-output.html](#)



Example 1

The Decision Boundary

Notebook: [decision-tree-real-input-real-output.html](#)



Objective Function for Regression Trees

Feature is denoted by X and target by Y .

Let the split be at $X = s$.

Define regions: $R_1 = \{x : x \leq s\}$ and $R_2 = \{x : x > s\}$.

Objective Function for Regression Trees

Feature is denoted by X and target by Y .

Let the split be at $X = s$.

Define regions: $R_1 = \{x : x \leq s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1} y_i$$

$$c_2 = \frac{1}{|R_2|} \sum_{x_i \in R_2} y_i$$

Objective Function for Regression Trees

Feature is denoted by X and target by Y .

Let the split be at $X = s$.

Define regions: $R_1 = \{x : x \leq s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1} y_i$$

$$c_2 = \frac{1}{|R_2|} \sum_{x_i \in R_2} y_i$$

The loss function is:

$$\text{Loss}(s) = \sum_{x_i \in R_1} (y_i - c_1)^2 + \sum_{x_i \in R_2} (y_i - c_2)^2$$

Objective Function for Regression Trees

Feature is denoted by X and target by Y .

Let the split be at $X = s$.

Define regions: $R_1 = \{x : x \leq s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1} y_i$$

$$c_2 = \frac{1}{|R_2|} \sum_{x_i \in R_2} y_i$$

The loss function is:

$$\text{Loss}(s) = \sum_{x_i \in R_1} (y_i - c_1)^2 + \sum_{x_i \in R_2} (y_i - c_2)^2$$

Our objective is to find the optimal split:

Algorithm: Finding the Optimal Split

1. Sort all data points (x_i, y_i) in increasing order of x_i .

Algorithm: Finding the Optimal Split

1. Sort all data points (x_i, y_i) in increasing order of x_i .
2. Evaluate the loss function for all candidate splits:

$$s = \frac{x_i + x_{i+1}}{2} \text{ for } i = 1, 2, \dots, n - 1$$

3. Select the split s^* that minimizes the loss function.

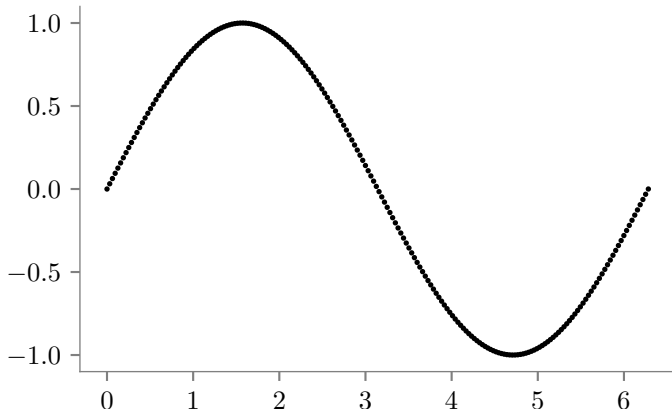
A Question!

Draw a regression tree for $Y = \sin(X)$, $0 \leq X \leq 2\pi$

A Question!

Dataset of $Y = \sin(X)$, $0 \leq X \leq 7$ with 10,000 points

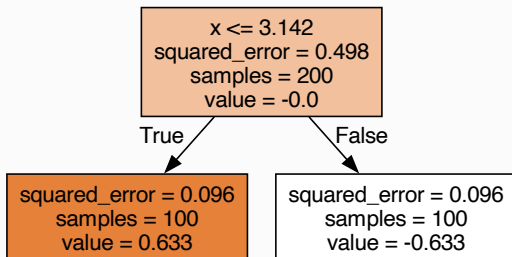
Notebook: [decision-tree-real-input-real-output.html](#)



A Question!

Regression tree of depth 1

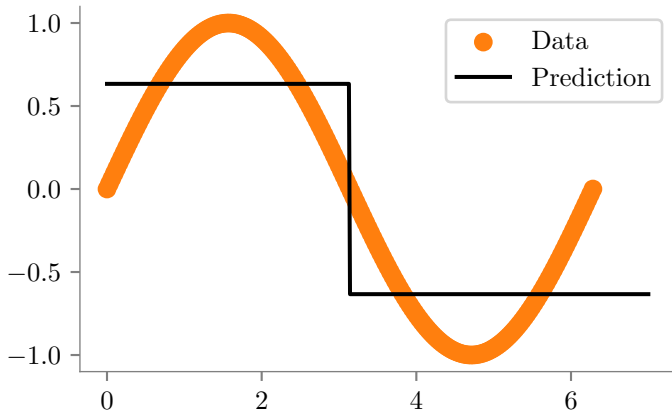
Notebook: decision-tree-real-input-real-output.html



A Question!

Decision Boundary

Notebook: [decision-tree-real-input-real-output.html](#)

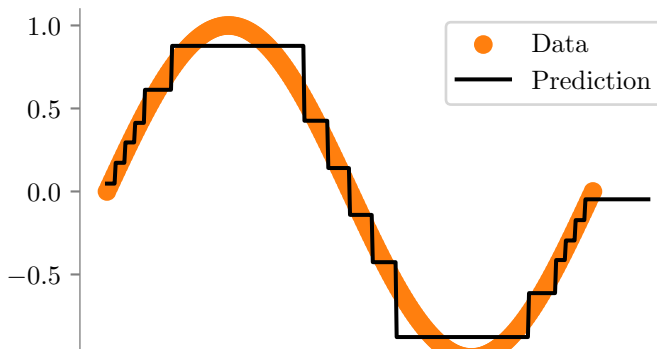


A Question!

Regression tree with no depth limit is too big to fit in a slide.

It has of depth 4. The decision boundaries are in figure below.

Notebook: [decision-tree-real-input-real-output.html](#)



Pop Quiz #31

Quick Question!

What is the prediction function for a regression tree leaf node?

A) The median of target values in that region

Pop Quiz #32

Quick Question!

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region

Pop Quiz #33

Quick Question!

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region

Pop Quiz #34

Quick Question!

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

Pop Quiz #35

Quick Question!

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

Pop Quiz #36

Quick Question!

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

Answer: C) The mean of target values in that region - Each leaf predicts the average target value of training samples that reach that leaf.

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data
- **Symptoms:**

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data
- **Symptoms:**
 - High training accuracy, low test accuracy

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data
- **Symptoms:**
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data
- **Symptoms:**
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves
 - Rules that are too specific to training data

The Problem: Overfitting in Decision Trees

- **Unpruned trees:** Can grow very deep and complex
- **Perfect training accuracy:** Each leaf contains single training example
- **But:** Poor generalization to new data
- **Symptoms:**
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves
 - Rules that are too specific to training data
- **Solution:** Pruning to control model complexity

Pre-pruning (Early Stopping)

Stop growing tree before it becomes too complex:

- **Maximum depth:** Limit tree depth (e.g., max_depth = 5)

Advantages: Simple, computationally efficient

Disadvantages: May stop too early, miss good splits later

Pre-pruning (Early Stopping)

Stop growing tree before it becomes too complex:

- **Maximum depth:** Limit tree depth (e.g., `max_depth = 5`)
- **Minimum samples per split:** Don't split if node has $< N$ samples

Advantages: Simple, computationally efficient

Disadvantages: May stop too early, miss good splits later

Pre-pruning (Early Stopping)

Stop growing tree before it becomes too complex:

- **Maximum depth:** Limit tree depth (e.g., `max_depth = 5`)
- **Minimum samples per split:** Don't split if node has $< N$ samples
- **Minimum samples per leaf:** Ensure each leaf has $\geq M$ samples

Advantages: Simple, computationally efficient

Disadvantages: May stop too early, miss good splits later

Pre-pruning (Early Stopping)

Stop growing tree before it becomes too complex:

- **Maximum depth:** Limit tree depth (e.g., `max_depth = 5`)
- **Minimum samples per split:** Don't split if node has $< N$ samples
- **Minimum samples per leaf:** Ensure each leaf has $\geq M$ samples
- **Maximum features:** Consider only subset of features at each split

Advantages: Simple, computationally efficient

Disadvantages: May stop too early, miss good splits later

Pre-pruning (Early Stopping)

Stop growing tree before it becomes too complex:

- **Maximum depth:** Limit tree depth (e.g., `max_depth = 5`)
- **Minimum samples per split:** Don't split if node has $< N$ samples
- **Minimum samples per leaf:** Ensure each leaf has $\geq M$ samples
- **Maximum features:** Consider only subset of features at each split
- **Minimum impurity decrease:** Only split if improvement $>$ threshold

Advantages: Simple, computationally efficient

Disadvantages: May stop too early, miss good splits later

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**
 1. Grow complete tree on training data

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**

1. Grow complete tree on training data
2. Use validation set to evaluate subtree performance

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**

1. Grow complete tree on training data
2. Use validation set to evaluate subtree performance
3. Remove branches that don't improve validation accuracy

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**

1. Grow complete tree on training data
2. Use validation set to evaluate subtree performance
3. Remove branches that don't improve validation accuracy
4. Repeat until no beneficial removals remain

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**

1. Grow complete tree on training data
2. Use validation set to evaluate subtree performance
3. Remove branches that don't improve validation accuracy
4. Repeat until no beneficial removals remain

- **Cost Complexity Pruning:** Minimize

Error + $\alpha \times \text{Tree Size}$

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**
 1. Grow complete tree on training data
 2. Use validation set to evaluate subtree performance
 3. Remove branches that don't improve validation accuracy
 4. Repeat until no beneficial removals remain
- **Cost Complexity Pruning:** Minimize
 $\text{Error} + \alpha \times \text{Tree Size}$
- **Advantages:** More thorough, can recover from early stopping mistakes

Post-pruning (Tree Simplification)

Grow full tree, then remove unnecessary branches:

- **Algorithm:**
 1. Grow complete tree on training data
 2. Use validation set to evaluate subtree performance
 3. Remove branches that don't improve validation accuracy
 4. Repeat until no beneficial removals remain
- **Cost Complexity Pruning:** Minimize
 $\text{Error} + \alpha \times \text{Tree Size}$
- **Advantages:** More thorough, can recover from early stopping mistakes
- **Disadvantages:** More computationally expensive

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- **Process:**

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- **Process:**
 1. Start with full tree ($\alpha = 0$)

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- **Process:**
 1. Start with full tree ($\alpha = 0$)
 2. Gradually increase α

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- **Process:**
 1. Start with full tree ($\alpha = 0$)
 2. Gradually increase α
 3. At each α , prune branches that increase cost

Cost Complexity Pruning Algorithm

Systematic approach to find optimal tree size:

- **Cost function:** $R_{\alpha}(T) = R(T) + \alpha|T|$
 - $R(T)$: Misclassification error on validation set
 - $|T|$: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- **Process:**
 1. Start with full tree ($\alpha = 0$)
 2. Gradually increase α
 3. At each α , prune branches that increase cost
 4. Select α with best cross-validation performance

Bias-Variance Trade-off in Trees

- **Unpruned trees:**

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**
 - High bias (may miss important patterns)

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting
- **Optimal pruning:** Balances bias and variance

Bias-Variance Trade-off in Trees

- **Unpruned trees:**
 - Low bias (can fit complex patterns)
 - High variance (sensitive to training data changes)
 - Prone to overfitting
- **Heavily pruned trees:**
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting
- **Optimal pruning:** Balances bias and variance
- **Cross-validation:** Essential for finding this balance

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters** (sklearn):

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters** (sklearn):
 - `max_depth`: Start with 3-10

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters (sklearn):**
 - `max_depth`: Start with 3-10
 - `min_samples_split`: Try 10-100

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters (sklearn):**
 - `max_depth`: Start with 3-10
 - `min_samples_split`: Try 10-100
 - `min_samples_leaf`: Try 5-50

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters (sklearn):**
 - `max_depth`: Start with 3-10
 - `min_samples_split`: Try 10-100
 - `min_samples_leaf`: Try 5-50
 - `ccp_alpha`: Use for cost complexity pruning

Practical Pruning Guidelines

- **Start simple:** Begin with restrictive pre-pruning parameters
- **Cross-validation:** Always use CV to select pruning parameters
- **Validation curves:** Plot training/validation error vs. tree complexity
- **Common parameters** (sklearn):
 - `max_depth`: Start with 3-10
 - `min_samples_split`: Try 10-100
 - `min_samples_leaf`: Try 5-50
 - `ccp_alpha`: Use for cost complexity pruning
- **Domain knowledge:** Consider interpretability requirements

Summary

- Interpretability an important goal

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize “performance gain”

Summary

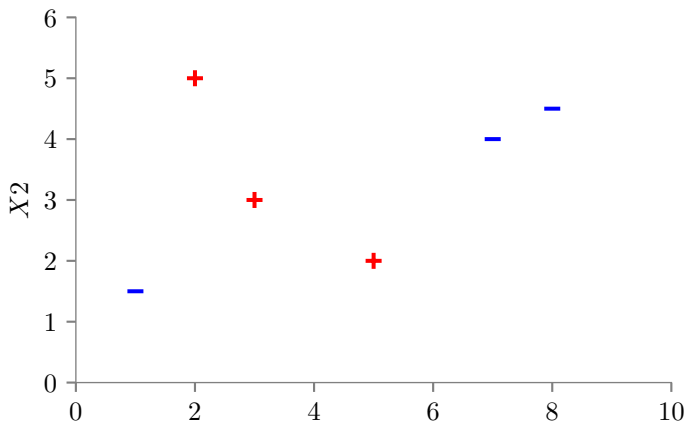
- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize “performance gain”
- Issues:

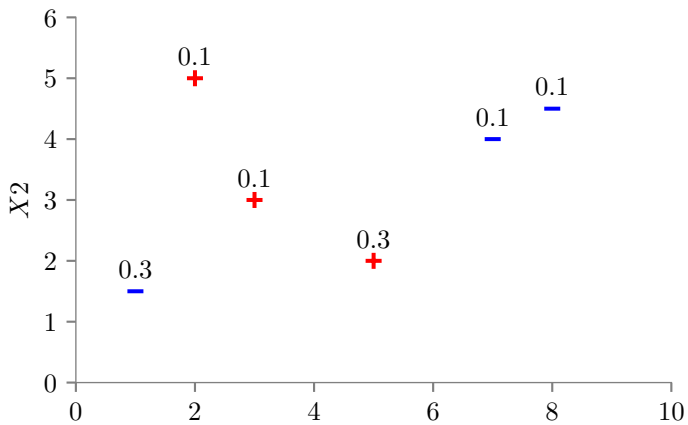
Summary

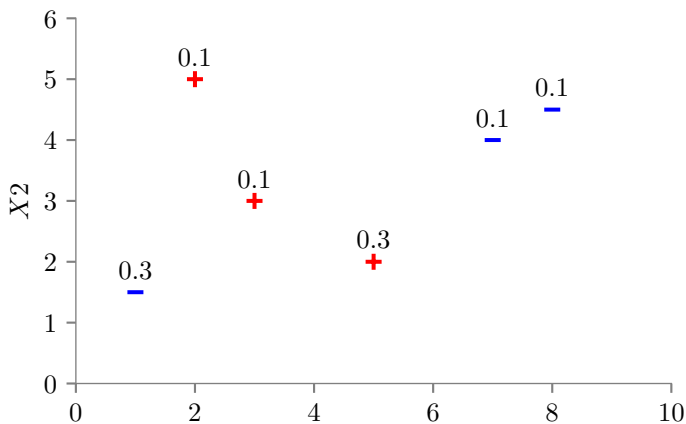
- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize “performance gain”
- Issues:
 - Can overfit easily!

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize “performance gain”
- Issues:
 - Can overfit easily!
 - Empirically not as powerful as other methods

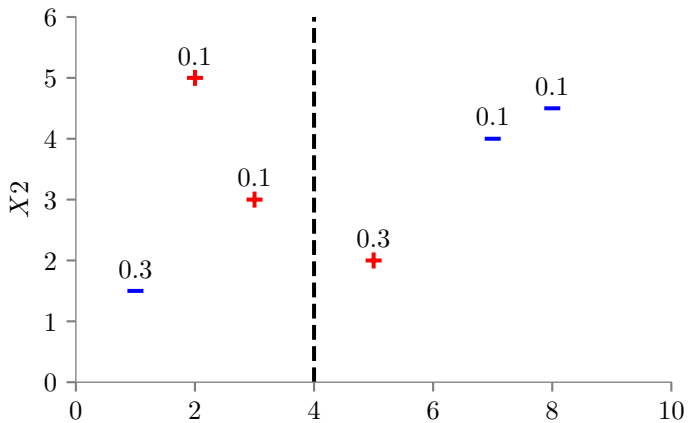




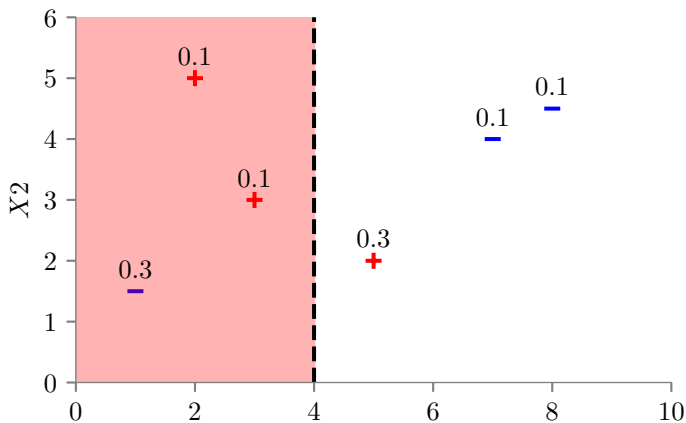


$$\text{Entropy} = -P(+) \log_2 P(+) - P(-) \log_2 P(-)$$

$$P(+) = \frac{0.1 + 0.1 + 0.3}{1} = 0.5, \quad P(-) = \frac{0.3 + 0.1 + 0.1}{1} = 0.5$$



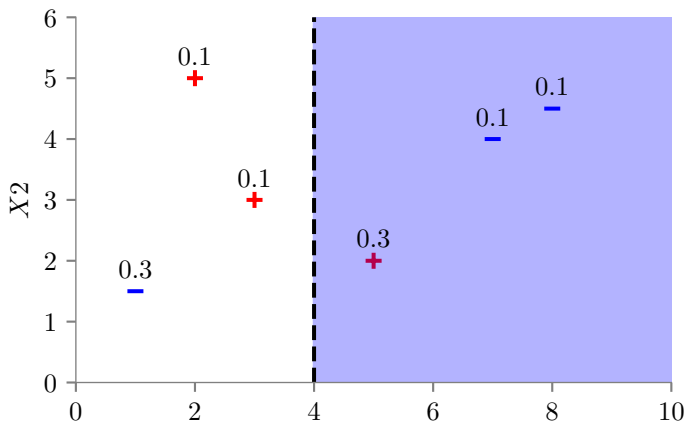
Candidate Line: $X_1 = 4(X_1^*)$



Entropy of $X_1 \leq X_1^* = E_{S(X_1 < X_1^*)}$

$$P(+)=\frac{0.1+0.1}{0.1+0.1+0.3}=\frac{2}{5}$$

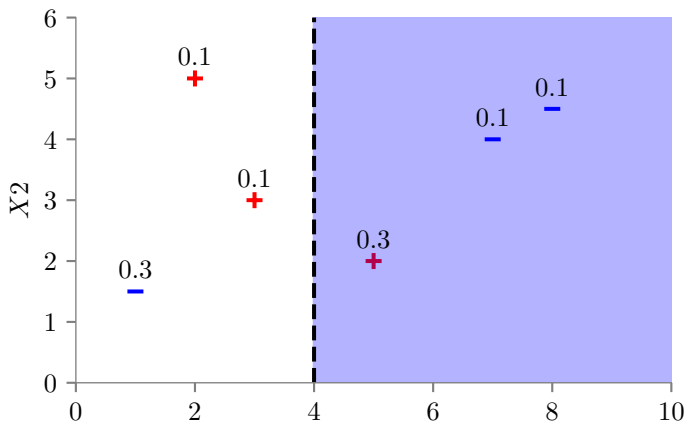
$$P(-)=\frac{3}{5}$$



Entropy of $X_1 > X_1^* = E_{S(X_1 > X_1^*)}$

$$P(+)=\frac{3}{5}$$

$$P(-)=\frac{2}{5}$$



$$\text{IG}(X_1 = X_1^*) = E_S - \frac{0.5}{1} \cdot E_{S(X_1 < X_1^*)} - \frac{0.5}{1} \cdot E_{S(X_1 > X_1^*)}$$